

1 **CLAIMS**

2 1. A method for deriving server resource utilization estimates for a
3 server cluster, the method comprising:

4 recording server cluster data during operation of the server cluster, at least
5 some of the server cluster data indicating server resource parameter values;

6 using a load simulation tool that, using the recorded data, determines a
7 maximum load that can be handled by the server cluster;

8 specifying a load to be handled by the server cluster; and

9 deriving server resource utilization estimates corresponding to the specified
10 load.

11
12 2. The method as recited in claim 1, further comprising:

13 displaying the server resource utilization estimates; and

14 recommending a plan to optimize processing of the specified load.

15
16 3. The method as recited in claim 2, wherein the plan recommends a
17 change in the hardware configuration of the server cluster.

18
19 4. The method as recited in claim 1, wherein the maximum load, the
20 recorded values, the specified load, and the server resource utilization estimates
21 are stored in non-volatile memory.
22
23
24
25

1
2 5. The method as recited in claim 1, wherein the using a load simulation
3 tool comprises:

4 creating a test script from the recorded values;
5 running the test script on a master client to simulate load and server
6 resource utilization conditions that existed on a server when the recorded values
7 were recorded; and
8 increasing the load on the server, when the test script is running, until a
9 maximum load that can be handled by the server is obtained.
10

11 6. The method as recited in claim 5, wherein the increasing the load on
12 the server further comprises multiplying the number of users utilizing the server
13 cluster when the recorded values were recorded, thereby multiplying the resources
14 utilized by the users.
15

16 7. The method as recited in claim 5, wherein the increasing the load on
17 the server further comprises decreasing the amount of time between user requests,
18 thereby increasing the resources utilized by the users.
19
20
21
22
23
24
25

1 **8.** The method as recited in claim 5, wherein the increasing the load on
2 the server until a maximum load that can be handled by the server is obtained,
3 further comprises:

4 observing a service rate exhibited by the server on which the simulation is
5 being performed; and

6 recognizing that the maximum load has been obtained when an increase in
7 the load does not increase the service rate.

8
9 **9.** The method as recited in claim 5, wherein:

10 the server cluster contains a set of identical servers;

11 running the test script run on the master client simulates server cluster
12 operation on only one of the servers of the server cluster; and

13 the method further comprises extrapolating the results obtained on the one
14 server using the number of servers in the set of identical servers to obtain the
15 maximum load that can be handled by the server cluster.

16
17 **10.** The method as recited in claim 5, wherein:

18 the server cluster contains a set of non-identical servers;

19 running the test script run on the master client further comprises running
20 the test script on each of the non-identical servers in the server cluster; and

21 the method further comprises summing the results obtained from each non-
22 identical server in the cluster to obtain the maximum load that can be handled by
23 the server cluster.

1 11. The method as recited in claim 1, wherein the recording server
2 cluster data during operation of the server cluster comprises recording data
3 directly from the server cluster and recording data that is input by a server cluster
4 user.

5
6 12. The method as recited in claim 1, wherein the server resource
7 utilization estimates comprise estimates for one or more of the following:
8 processor utilization, memory utilization, communication bandwidth utilization,
9 and general server utilization.

10
11 13. The method as recited in claim 1, wherein the server resource
12 utilization comprises general server utilization, the method further comprising:

13 Deriving general server utilization by solving:

14
15
$$U = \frac{L}{X}$$

16

17
18 Wherein U is the general server utilization; X is the maximum load that can
19 be handled by the server cluster; and L is the specified load.
20
21
22
23
24
25

1
2 14. The method as recited in claim 1, wherein the server resource
3 utilization comprises processor utilization, the method further comprising:

4 deriving processor utilization by solving:

5
6
$$U_{CPU} = \frac{a}{e^{b \cdot L}}$$

7
8

9 wherein U_{CPU} is processor utilization; L is the specified load; a is processor
10 regression constant a ; and b is processor regression constant b .
11

12 15. The method as recited in claim 1, wherein the server resource
13 utilization comprises communication bandwidth utilization, the method further
14 comprising:

15 deriving communication bandwidth utilization by solving:

16
17
$$U_B = \frac{F_{TCP}}{B} \cdot (c + d \cdot L)$$

18
19

20 wherein U_B is communication bandwidth utilization; L is the specified load;
21 c is processor regression constant c ; d is processor regression constant d ; F_{TCP} is a
22 transmission overhead factor; and B is the total communication bandwidth
23 available.
24
25

1 16. The method as recited in claim 1, wherein the server resource
2 utilization comprises memory utilization, the method further comprising:

3
4 deriving memory utilization by solving:

5
6
$$U_M = \frac{N \cdot (M_{TCP} + M_{IISStruct}) + M_{OS} + M_{IIS}}{M}$$

7

8 wherein N is a total number of concurrent connections derived by solving:
9

10
$$N = \frac{L}{(X - L)} + S1 \cdot L$$

11

12 wherein: U_M is memory utilization; M_{TCP} is a an amount of memory
13 necessary to support the connections for communications; $M_{IISStruct}$ is the amount
14 of memory necessary to support data structures associated with each connection;
15 M_{OS} is the amount of memory required by a server operating system; M_{IIS} is the
16 amount of memory required by a server communication program; M is the total
17 amount of memory available; L is the specified load; X is the maximum load that
18 can be handled by the server cluster; and $S1$ is a connection memory factor that is
19 the adjusted average of the incoming connections at different speeds.
20
21
22
23
24
25

1
2 17. The method as recited in claim 1, wherein the server resource
3 utilization comprises general server utilization, the method further comprising:

4
5 18. The method as recited in claim 1, wherein the server resource
6 utilization comprises processor utilization, the method further comprising:

7 deriving processor utilization by solving:

8
9
10
$$U_{CPU} = \frac{a}{e^{b \cdot L}}$$

11
12 wherein U_{CPU} is processor utilization; L is the specified load; a is processor
13 regression constant a ; and b is processor regression constant b .
14
15
16
17
18
19
20
21
22
23
24
25

1
2 19. The method as recited in claim 1, wherein the server resource
3 utilization comprises communication bandwidth utilization, the method further
4 comprising:

5 deriving communication bandwidth utilization by solving:

6
7
$$U_B = \frac{F_{TCP}}{B} \cdot (c + d \cdot L)$$

8

9 wherein U_B is communication bandwidth utilization; L is the specified load;
10 c is processor regression constant c ; d is processor regression constant d ; F_{TCP} is a
11 transmission overhead factor; and B is the total communication bandwidth
12 available.
13
14
15
16
17
18
19
20
21
22
23
24
25

1
2 **20.** The method as recited in claim 1, wherein the server resource
3 utilization comprises memory utilization, the method further comprising:

4
5 deriving memory utilization by solving:

6
7
$$U_M = \frac{N \cdot (M_{TCP} + M_{IISStruct}) + M_{OS} + M_{IIS}}{M}$$

8

9
10 wherein N is a total number of concurrent connections derived by solving:

11
$$N = \frac{L}{(X - L)} + S1 \cdot L$$

12

13 wherein: U_M is memory utilization; M_{TCP} is a an amount of memory
14 necessary to support the connections for communications; $M_{IISStruct}$ is the amount
15 of memory necessary to support data structures associated with each connection;
16 M_{OS} is the amount of memory required by a server operating system; M_{IIS} is the
17 amount of memory required by a server communication program; M is the total
18 amount of memory available; L is the specified load; X is the maximum load that
19 can be handled by the server cluster; and $S1$ is a connection memory factor that is
20 the adjusted average of the incoming connections at different speeds.
21
22
23
24
25

1
2 **21.** The method recited in claim 1, wherein the server resource
3 utilization is processor utilization, the method further comprising:

4 finding a functional dependency approximation between processor
5 utilization and load;

6 transforming functional dependency into linear form by using logarithmic
7 transformation;

8 deriving first and second processor regression constants using linear
9 regression methodology;

10 dividing the first processor regression constant by e to the power of the
11 product of the second processor regression constant and the specified load to
12 obtain the processor utilization estimate.

13
14 **22.** The method as recited in claim 1, wherein the server resource
15 utilization is communication bandwidth utilization, the method further comprising:

16 finding a functional dependency approximation between communication
17 bandwidth utilization;

18 transforming functional dependency into linear form by using logarithmic
19 transformation;

20 deriving first and second bandwidth regression constants using linear
21 regression methodology;

22 deriving a transmission overhead factor that, when applied to a certain size
23 web page, results in the actual capacity necessary to transmit the web page;

24 deriving a weighted communication overhead factor by dividing the
25 transmission overhead factor by the available communication bandwidth;

1 deriving an adjusted communication load by adding the first bandwidth
2 regression constant to the product of the specified load and the second bandwidth
3 regression constant; and

4 determining the communication bandwidth utilization estimate by
5 multiplying the weighted communication overhead factor by the adjusted
6 communication load.

7
8 **23.** The method as recited in claim 1, wherein the server resource
9 utilization is memory utilization, the method further comprising:

10 deriving a connection memory factor that is the adjusted average of the
11 incoming connections at different speeds;

12 deriving a weighted connection memory factor by multiplying the
13 connection memory factor by the specified load;

14 deriving a page load ratio by dividing the specified load by the difference of
15 the maximum load value and the specified load;

16 deriving a total number of concurrent connections by adding the weighted
17 connection memory factor and the page load ratio; and

18 deriving a gross memory utilization by multiplying the total number of
19 concurrent connections by the sum of the amount of memory necessary to support
20 each connection for communications and the amount of memory necessary to
21 support data structures associated with each connection, and adding the amount of
22 memory required by a server operating system and the amount of memory
23 required by the server communication program; and

24 deriving the memory utilization estimate by dividing the gross memory
25 utilization by total memory available.

1
2 **24.** The method as recited in claim 1, wherein the server resource
3 utilization is general server utilization, the method further comprising:

4 dividing the specified load by the maximum load to derive the general
5 server utilization estimate.

6
7 **25.** One or more computer-readable media having computer-readable
8 instructions thereon which, when executed by one or more computers, cause the
9 computers to implement the method of claim 1.

10
11 **26.** A simulation tool for use in determining server resource utilization
12 estimates in a server cluster having one or more servers, the load simulation tool
13 comprising:

14 a user interface configured to receive data input from a user;

15 at least one filter or monitor configured to record operational data from one
16 or more of the servers in the server cluster;

17 the simulation tool being configured to create a test script from the recorded
18 data and the received data, and to run the test script from a master client connected
19 to the server cluster to simulate load and other server conditions that existed when
20 the operational data was recorded; and

21 the user interface being further configured to display utilization of server
22 resources during the running of the test script.

1 27. The simulation tool as recited in claim 26, wherein the simulation
2 tool being configured to create the test script is further configured to allow the user
3 to create the test script to increase the load on the server on which the simulation is
4 running and observe the effect of such an increase in load on the server resource
5 utilization displays.

6
7 28. A system, comprising:
8 a server cluster having one or more servers, one of which is a primary
9 server that controls the operation of the server cluster;
10 a cluster controller resident in memory on the primary server of the server
11 cluster, the cluster controller controlling communications between the primary
12 server and secondary servers, if any, and between clients and the server cluster;
13 an operating system resident in the memory of the primary server;
14 a communications program within the cluster controller to provide
15 communications capability for the system;
16 a filter to collect server data indicating certain operating parameters for the
17 server cluster;
18 a monitor on each server in the server cluster to collect server data
19 indicating certain operating parameters for the server cluster;
20 a user interface to collect data input by a user;
21 a capacity planner within the cluster controller configured to utilize the
22 collected data to derive one or more server resource utilization estimates for server
23 resources to determine how handling a specified load will affect the utilization of
24 the server resources, and to produce a plan recommending changes to be made to
25 the server cluster to adequately accommodate the specified load; and

1 a load simulation tool configured to use the collected data to create a
2 simulation script that, when run on a master client, simulates the operation of the
3 server cluster system to allow the user to find the maximum load that the server
4 cluster can handle; and

5 wherein the maximum load obtained through the use of the load simulation
6 tool is utilized in the derivation of the one or more server resource utilization
7 estimates.

8
9 **29.** The system as recited in claim 28, wherein the filter is an ISAPI
10 filter.

11
12 **30.** The system as recited in claim 28, wherein the collected data and the
13 plans are stored in the memory.

14
15 **31.** The system as recited in claim 28, wherein the simulation script is
16 run from a master client connected to the server cluster, and wherein the
17 simulation is performed on only one server of the server cluster.

18
19 **32.** The system as recited in claim 31, wherein the load simulation tool
20 is further configured to extrapolate results from the simulation on one server in the
21 server cluster to derive results for the total number of servers in the server cluster.
22
23
24
25

1
2 **33.** The system as recited in claim 28, wherein:
3 the load simulation tool is further configured to run a simulation script from
4 a master client connected to the server cluster;
5 the simulation is performed on each server in the server cluster; and
6 results from each server are summed to derive results of the total number of
7 servers in the server cluster.

8
9 **34.** The system as recited in claim 28, wherein the server resource
10 utilization derived by the capacity planner comprises general server utilization,
11 and the capacity planner is further configured to derive general server utilization
12 by solving:

13
14
$$U = \frac{L}{X}$$

15
16

17 wherein U is the general server utilization; X is the maximum load that can
18 be handled by the server cluster which is determined by the load simulation tool;
19 and L is the specified load.
20
21
22
23
24
25

1
2 35. The system as recited in claim 28, wherein the server resource
3 utilization derived by the capacity planner comprises general server utilization,
4 and the capacity planner is further configured to derive general server utilization
5 by solving:

6
7
8
$$U_{CPU} = \frac{a}{e^{b \cdot L}}$$

9

10 wherein U_{CPU} is processor utilization; L is the specified load; a is processor
11 regression constant a ; and b is processor regression constant b .
12

13 36. The system as recited in claim 28, wherein the server resource
14 utilization derived by the capacity planner comprises communication bandwidth
15 utilization, and the capacity planner is further configured to derive communication
16 bandwidth utilization by solving:
17

18
$$U_B = \frac{F_{TCP}}{B} \cdot (c + d \cdot L)$$

19
20

21 wherein U_B is communication bandwidth utilization; L is the specified load;
22 c is processor regression constant c ; d is processor regression constant d ; F_{TCP} is a
23 transmission overhead factor; and B is the total communication bandwidth
24 available.
25

1
2 37. The system as recited in claim 28, wherein the server resource
3 utilization derived by the capacity planner comprises communication bandwidth
4 utilization, and the capacity planner is further configured to derive communication
5 bandwidth utilization by solving:

6
7
$$U_M = \frac{N \cdot (M_{TCP} + M_{IISStruct}) + M_{OS} + M_{IIS}}{M}$$

8

9 wherein N is a total number of concurrent connections derived by solving:
10

11
12
$$N = \frac{L}{(X - L)} + S1 \cdot L$$

13

14
15 wherein: U_M is memory utilization; M_{TCP} is a an amount of memory
16 necessary to support the connections for communications; $M_{IISStruct}$ is the amount
17 of memory necessary to support data structures associated with each connection;
18 M_{OS} is the amount of memory required by a server operating system; M_{IIS} is the
19 amount of memory required by a server communication program; M is the total
20 amount of memory available; L is the specified load; X is the maximum load that
21 can be handled by the server cluster; and $S1$ is a connection memory factor that is
22 the adjusted average of the incoming connections at different speeds.
23
24
25